

1 **Widespread of horizontal gene transfer regions in eukaryotes**

2 Kun Li^{1§}, Fazhe Yan^{1§}, Zhongqu Duan¹, David L. Adelson², Chaochun Wei^{1, 3*}

3 ¹ School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, 800

4 Dongchuan Road, Shanghai 200240, China

5 ² School of Biological Sciences, The University of Adelaide, SA 5005, Australia

6 ³ Joint International Research Laboratory of Metabolic and Developmental Sciences,

7 Shanghai Jiao Tong University, Shanghai 200240, China

8 **Contact information**

9 Chaochun Wei

10 Shanghai Jiao Tong University, 800 Dongchuan Road, Shanghai 200240, China

11 Tel: (+86)21-34204083

12 E-mail: ccwei@sjtu.edu.cn

13 §: these authors contributed equally to this work.

14 **Summary**

15 Horizontal gene transfer (HGT) is the transfer of genetic material between distantly related
16 organisms. While most genes in prokaryotes can be horizontally transferred, HGTs in eukaryotes
17 are considered as rare, particularly in mammals. Here we report the identification of HGT regions
18 in 13 model eukaryotes by comparing their genomes with 824 eukaryotic genomes. Between 4 and
19 358 non-redundant HGT regions per species were found in the 13 model organisms, and most of
20 these HGT regions were previously unknown. The majority of the 824 eukaryotes with full length
21 genome sequences also contain HGTs. These HGTs have transformed their host genomes with
22 thousands of copies and have impacted hundreds, even thousands of genes. We extended this
23 analysis to ~128,000 prokaryote and virus genomes and revealed a few potential routes of horizontal
24 gene transfer involving blood sucking parasites, intracellular pathogens, and bacteria. Our findings
25 revealed that HGT is widespread in eukaryotic genomes, and is an important driver of genome
26 evolution for eukaryotes.

27

28 **Main**

29 Horizontal gene transfer (HGT) is the transfer of genetic material between organisms that is not
30 from parent to offspring, and it is a major driver of genome evolution in bacteria and archaea^{1,2}. On
31 average, 81% of genes in prokaryotes were involved in HGT³. Recent evidence has shown that
32 HGTs also exist in eukaryotes. For example, HGTs have been reported from soil bacteria to the
33 common ancestor of *Zygnematophyceae* and *embryophytes*, which increased its resistance to biotic
34 and abiotic stresses during terrestrial adaptation⁴. Besides, HGT of a plant detoxification gene
35 BtPMaT1, made whiteflies gain the ability to malonylate a common group of plant defense
36 compounds⁵. Another remarkable example of HGT is a ~1.5Mb fragment of *Wolbachia spp.* DNA
37 integrated into the pill bug *Armadillidium vulgare* genome, resulting in the creation of a new W sex
38 chromosome⁶. HGTs have been observed in five parasitic plants in the *Orobanchaceae* family⁷,
39 several unicellular pathogens⁸ and blood-sucking parasites⁹⁻¹¹. Although it has been proposed that
40 only unicellular and early developmental stages in eukaryotes are vulnerable to HGT¹², some argue
41 that HGTs in eukaryotes may be limited to those derived from endosymbiotic organelles^{13,14}.

42 Mechanisms for the transfer of DNA into eukaryotic genomes have been described for viral
43 infection, transposons, conjugation between bacteria and eukaryotes, or from endosymbionts (not
44 only plastids and mitochondria)¹⁵. Some behaviors, such as predation, and life-styles, such as
45 parasitism, have been reported to promote DNA transfer in eukaryotes¹⁶. Recently, more eukaryote
46 HGTs were reported. For example, a plant detoxification gene BtPMT1 was found transferred to
47 whitefly *Bemisia tabaci* which greatly expanded the insect's food spectrum¹⁷. Therefore, while the
48 prevalence of HGT may be rare in eukaryotes, compared to bacteria and archaea, it does occur.
49 However, the scale and impact of HGTs in eukaryotes are unknown.

50 We present here a fast identification method for HGTs in eukaryotes using both sequence
51 composition bias and genome comparisons (Fig S1; see Methods), and we evaluated the method
52 using a simulated dataset. We applied this to 13 model organisms with high quality genomes, and
53 then expanded it to 824 eukaryotes with available full-length genome sequences. Many bacteria and
54 virus genomes were also compared.

55

56 **Results**

57 **A Fast HGT identification method and evaluation of the method.**

58 We created a fast identification pipeline for HGTs in eukaryotes by combining sequence
59 composition filtering and genome sequence comparison (Fig S1; see Methods). The pipeline was
60 evaluated HGTs previously reported HGTs in the human genome¹⁹. We tried different k-mer sizes
61 (1~6), and k=4 was selected because the highest portion of candidate HGTs previously reported in
62 the human genome¹⁹ were kept (Figure S6). A very high portion (>75%) of these human HGTs
63 reported previously were kept in the result HGTs even if we only input top 5% of the fragments
64 with highest differences to the human genome.

65

66 We further evaluated the pipeline with a genome containing simulated HGTs. Since our HGT
67 identification pipeline has two main steps, sequence composition-based filtering step and genome
68 comparison step. The evaluation was done for the two steps (Figure S8). While top 1% fragments
69 were input to the pipeline, 20.6% correct results would be identified after sequence composition-

70 based filtering and 14.3% correct results identified after genome comparison. When the percentage
71 of fragments input was up to 50%, 83.4% and 77.7% correct results were identified after two steps
72 respectively. It can be seen that the prediction precisions were higher than 60% in all cases. This
73 indicated that we may have underestimated the number of HGTs (low recall rate) but the identified
74 HGTs were highly reliable.

75

76 **Widespread of HGTs among eukaryotes**

77 We applied our HGT identification method to identify HGTs in Eukaryotes. For the 13 model
78 organisms with high quality genomes ([Table 1](#)), we identified between 4 and 358 non-redundant
79 eukaryotic HGTs. For 824 eukaryotes with full length genome sequences currently available, almost
80 all (98.7%) of them contained HGTs. A number of those HGTs were also found to have bacteria or
81 viruses as media vectors ([Table 1](#), [Table S1](#)).

82 For the identified HGTs in the 13 model organisms ([Table 1](#)), most of them were previously
83 unknown compared with reported HGTs^{10,18-32}. For each candidate HGT region, a phylogenetic tree
84 was constructed from the homologous sequences of that HGT region in all eukaryotes ([see](#)
85 [Methods](#)). To determine the frequency of HGT in eukaryotes, we calculated an HGT-appearance
86 number N_{AB} for a model organism A and another eukaryotic organism B, which was defined by the
87 frequency with which organism B appeared in the phylogenetic trees of non-redundant HGTs of
88 model organism A. For instance, among the 313 non-redundant HGT trees for *Homo sapiens*, *Pan*
89 *trogodytes* was found in 312 of them, therefore the HGT-appearance number N_{HP} between *Homo*
90 *sapiens* and *Pan troglodytes* was 312. The distribution of HGT-appearance numbers between the 13
91 model organisms and 824 eukaryotes is shown in [Figure 1A](#) and [Table S2](#). If model organism A and
92 organism B were from different kingdoms, N_{AB} are shown as a line in [Figure 1B](#) and [Table S3](#). The
93 greater the value for N_{AB} , the thicker the line. By using this metric, we determined that 98.7% of
94 eukaryotes (813 of 824) hosted HGTs, revealing widespread of HGTs across eukaryotes. In
95 addition, we categorized the HGTs into cross-kingdom, cross-phylum, cross-class or unknown
96 categories based on the taxonomy relationships of the two involved organisms. We found that 1081
97 pairs of cross-kingdom species contained at least one HGT, and about half of them contained
98 multiple HGTs ([Figure 1B](#), [Table S3](#)). The number of cross-phylum and cross-class species pairs

99 containing HGT were 1,890 and 2,909 respectively ([Figure S2, Table S3](#)).

100 **Duplications of HGTs and their impact on their host genomes**

101 Horizontal transferred active transposable elements may continue to transpose in the new host.
102 Therefore, we compared the non-redundant eukaryotic HGT sequences we identified with their host
103 genomes. Overall, about 22.2% of HGTs (242 of 1,090) have multiple copies in their host genomes,
104 and 47 HGTs have more than 1,000 copies ([Table S1](#)). In particular, BovB related HGT region
105 “chr8:96500648-96500854” in the cow genome has 56,890 copies (total length 11.3 Mbp), which
106 is consistent with a previous study that BovB are present as many copies⁹. In newly identified HGTs,
107 elephant HGT region “scaffold_90:4162401-4162729” has 13,484 copies and occupies 0.15% of
108 the elephant genome (4.4Mbp/3.2Gb). Frog HGT region “chr1:8559133-8559400” has 7,027 copies
109 and occupies 1.7Mb.

110 These HGT copies have affected many genes as well. There are 51 HGT regions, each of which
111 impacts more than 100 protein coding genes in their host genome ([Table S1](#)). For example, the frog
112 HGT region mentioned above and its copies overlap (with at least 1bp) more than 10% of all protein-
113 coding genes (2,149 of 19,983), which is a dramatic impact on the frog genome functions. HGTs
114 with similar (but different degree of) impact on genome functions can also be found for most of the
115 13 model organisms. Especially cow, human, frog, elephant, zebrafish, and lizard, each of them has
116 more than 100 genes impacted by HGTs. More information can be found in [Table S1](#).

117 **Repetitive sequence composition of HGTs**

118 We compared the non-redundant HGTs detected in 13 model organisms with the repetitive
119 sequences annotated in their reference genomes. Between 0~100% of their HGTs overlapped with
120 interspersed repeats (excluding simple repeats) ([Table 1](#)), revealing significant species and repeat-
121 specificity ([Fig S3A](#)). The types of repeats overlapping with HGTs showed significant correlation
122 with overall genomic repeat composition ([Fig S3B](#)). Retrotransposons (SINEs and LINEs) were
123 common in HGTs detected in mammals, consistent with their frequencies in their host genomes. In
124 a frog genome (*Xenopus tropicalis*), DNA transposons, the main repeats for that genome, were
125 frequently found in HGTs. In comparison, in the rat genome (*Rattus norvegicus*), the distribution of
126 DNA transposons in HGTs was not consistent with their distribution in the host genomes. In the rat
127 genome, DNA transposons appeared in as many as 6 non-redundant HGTs (46%), while that repeat

128 only accounted for 3.1% of repeats in the genome.

129 BovB and L1 retrotransposons are the two most abundant transposable elements (TEs) in
130 ruminant and afrotherian genomes and replicate via an RNA intermediate³³. The horizontal transfer
131 of BovB is known to be widespread in animals¹⁰ and horizontal transfer of L1 has been shown in
132 plants, animals and several fungi¹⁸. In total, 44 of our non-redundant HGTs overlapped with BovB
133 retrotransposons in *Bos taurus* (Ruminantia) and *Loxodonta africana* (Afrotheria) (Table 1),
134 supporting previous results for horizontal transfer of BovB^{10,18}. Furthermore, 461 L1 horizontal
135 transfer events were identified in five mammals (Cow, Human, Elephant, mouse, and rat), providing
136 more evidence that L1 elements are horizontally transferred¹⁸. Surprisingly, 95.2% (20 of 21) (Table
137 S4) HGTs that overlapped with BovB retrotransposons in *Bos taurus*, were associated with the
138 possible intermediary species, the blood-sucking parasite *Cimex lectularius* (bed bug), which has
139 been reported by Ivancevic et al.¹¹. *Cimex lectularius* is known to feed on animal blood and can host
140 over 40 zoonotic pathogens³⁴, thus transmitting many infectious diseases³⁵. Figure 2A showed the
141 tree from bovine HGT region “chr25:1343971-1344200” and its homologs. In addition to the
142 candidate vector species, this HGT tree also included 12 mammals (9 Ruminantia, 2 Metatheria and
143 *Macaca mulatta*) and 5 non-mammalian vertebrates (2 fishes, 3 reptiles), which were clearly
144 clustered in distinct branches (Table S5). In addition, we identified these mobile DNA sequences in
145 several bacteria, including *Enterococcus faecium*, *Mycolicibacterium malmesburyense*, *Escherichia*
146 *coli* and *Anaplasma phagocytophilum* (Table S6). Using WGS data, we confirmed high similarity
147 homologs (sequence coverage>80%, sequence identity>90%) of this HGT region from *Cimex*
148 *lectularius* (NW_014465023.1|11681076-11681736) in 10 samples (collected from PRJNA259363,
149 PRJNA167477 and PRJNA432971, sequenced in different centers) (Table S7). Like in other bugs,
150 it appears that *Cimex lectularius* transferred DNA between the hosts it feeds on.

151 **Apicomplexan intracellular pathogens often participate in HGTs**

152 A considerable number of genes of intracellular pathogens have been acquired through HGT,
153 including Apicomplexa^{8,36}. In particular, *Toxoplasma gondii* is an obligate intracellular,
154 apicomplexan parasite that causes the disease toxoplasmosis in a wide range of warm-blooded
155 animals including humans^{37,38}, where it has been reported to infect up to one third of the world’s
156 population³⁹. About 0.21% of *Toxoplasma gondii* protein-coding genes were acquired through

157 HGT⁴⁰. Our analysis identified 401 HGTs from 11 model organisms and 25 apicomplexans (Table
158 1). *Toxoplasma gondii* ME49, *Plasmodium vivax* and *Plasmodium knowlesi* strain H appeared more
159 frequently in HGTs (Fig 3A).

160 *Toxoplasma gondii* ME49 participated in 218 human HGTs correlated with Apicomplexan
161 intracellular pathogens, making these cross-kingdom HGTs (Fig 3B). For instance, the HGT tree of
162 HGT region “chr11:24184801-24185043” shown in Figure 2B includes 1 Apicomplexan pathogen,
163 2 invertebrates and 2 non-mammalian vertebrates (Table S5). This HGT tree is inconsistent with
164 the phylogenetic tree of these organisms, and this HGT was also found in 36 bacterial strains,
165 indicating that these same DNA sequences were able to jump into bacteria as well as eukaryotes.
166 The apicomplexan pathogen (*Toxoplasma gondii* ME49) and a blood-sucking parasite (*Ixodes*
167 *scapularis*) are good candidate sources/vectors for DNA transfer into the human genome. Several
168 primates including *Homo sapiens*, *Gorilla gorilla gorilla*, *Pan troglodytes*, *Pongo abelii* and *Pan*
169 *paniscus* were clustered into a branch in the HGT tree, indicating that this DNA transfer event
170 happened in their common ancestor. Using WGS data, we successfully confirmed homologous
171 sequences in *Toxoplasma gondii* ME49, which further supported this DNA transfer event.

172

173 Discussion

174 HGTs are widespread in eukaryotes (in the 13 model organisms we examined in this study and
175 98.7% of other eukaryotes with whole genome sequences). Compared to HGTs in prokaryotes, the
176 number of non-redundant eukaryotic HGTs (4~358 regions) detected in these model organisms was
177 very small. In addition, we found many HGT regions by comparing a small part of the genome
178 sequences that were significantly different from their reference genomes. It is conceivable that the
179 number of HGT regions is much larger than this.

180 As shown in Figure 1A, the HGT-appearance numbers decreased when the phylogenetic
181 distance between two organisms increased. For example, primates appeared in most HGT trees for
182 human, followed by mammals. Most primates appeared in almost all HGT trees of *Homo sapiens*,
183 indicating that most these HGT sequences were inserted into the genome of their common ancestor.
184 We observed a similar distribution of HGT-appearance numbers in other model organisms,

185 indicating that most HGT regions identified by our pipeline were transferred before the divergence
186 of model organisms and their sibling lineages. This also implied that these HGTs may have
187 important functions as they have persisted¹.

188 For mammals (human, mouse, rat, cow, and elephant), we investigated the geographic
189 distribution of the two organisms involved for each HGT event. Most of the organism pairs were
190 from the same continent. For the 259 species related to HGTs that occurred to the common ancestor
191 of mammals, 213 (82.2%) species were located in the same continent as the corresponding model
192 organism, 43 (16.6%) species were not, and 3 (1.2%) species were undetermined (Table S8). The
193 continents began to separate about 200 Mya, around the same time that the oldest mammals
194 emerged^{41,42}. For 371 species related to HGTs that occurred into the common ancestors of the orders
195 of the model organisms, 357(96.2%) of them were found in the same continent with the
196 corresponding model organisms and 14 (3.8%) species were not, which were all related to HGTs of
197 the elephant (Table S8). Proboscidea, the order to which elephants belong, originated 55 Mya⁴³,
198 significantly later than the time that the continents separated.

199 Our study uncovered several putative routes for the exchange of genetic material between
200 distantly related eukaryotes. We propose that blood-sucking parasites (like *Cimex lectularius* and
201 *Ixodes scapularis*) and intracellular pathogens (like *Toxoplasma gondii* ME49) were involved in
202 DNA transfer between mammals and other eukaryotes and these transferred DNA sequences were
203 also found in pathogenic bacteria, suggesting exchange of genetic material between eukaryotes and
204 bacteria (Table S9). In this fashion, bacteria might serve as the vector for DNA transfer between
205 distantly related and eukaryotes that might not be in close contact with each other. We also found
206 highly similar homologous sequences in viral genomes for three HGTs in human, indicating that
207 viruses might be agents for integration of transferred DNAs in to eukaryotic genomes. Taken
208 together, these findings revealed a putative route for DNA transfer between distantly related
209 eukaryotes (Figure 4). Nevertheless, we observed that about 54.4% of HGTs events could be
210 interpreted by bacteria medium (Table 1). However, the detailed routes for DNA transfer for the
211 majority of HGT regions in this report are still unclear. With the progress of sequencing technology,
212 especially third generation sequencing technologies, high quality whole genome sequences can be
213 obtained for several HGT related species distributed across the tree of life, and this will provide a

214 good opportunity to determine the route and direction of HGT .

215 Functional annotation for genes overlapping with HGTs (see Methods) revealed some
216 significantly enriched Gene Ontology terms (GO terms) (Bonferroni<0.05) for protein-coding genes
217 from mouse, fruit fly and nematode as well as non-coding genes from yeast. (Table S10). The
218 significant GO terms for nematode were “hemidesmosome, intermediate filament”, while the
219 significant GO term for mouse was “protein kinase A binding”. HGTs in fruit fly that overlapped
220 with coding genes were enriched for “ATP binding, lipid particle, microtubule associated complex”,
221 etc. HGTs in yeast overlapped with non-coding genes enriched for “retrotransposon nucleocapsid,
222 transposition, RNA-mediated, cytosolic large ribosomal subunit”, etc.

223 In conclusion, comparison of 13 model eukaryote genomes against other organisms with whole
224 available genome sequences showed that HGT is widespread in eukaryotes. We suggest that
225 blood-sucking parasites, apicomplexan pathogens, bacteria, and viruses are nodes in the putative
226 routes for DNA transfer between distantly related eukaryotes.

227

228 **Tables**

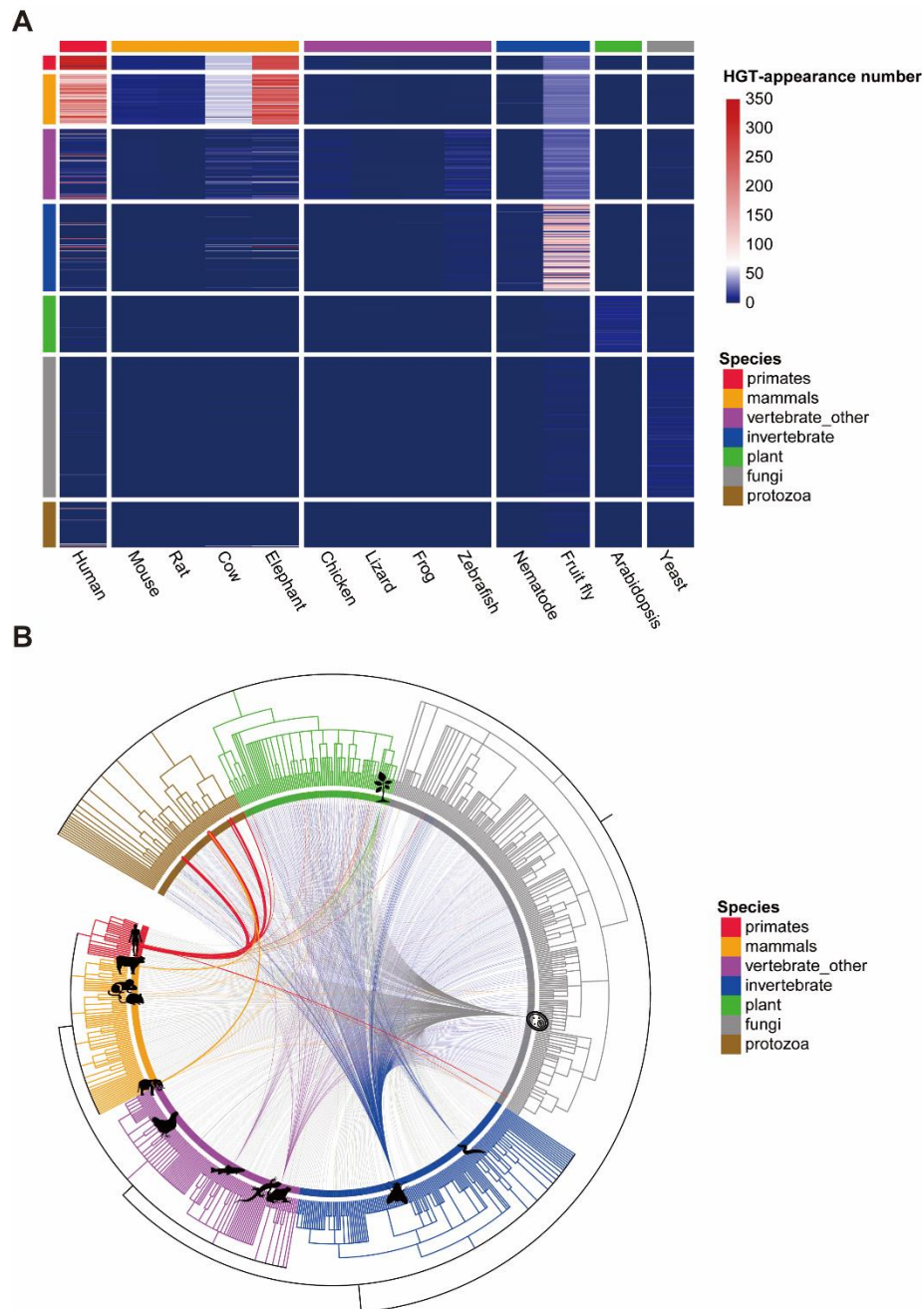
229 **Table 1. The numbers of non-redundant HGTs in 13 model organisms.** Most of these HGTs
 230 were novel. Some of these HGTs are supported by genomic evidence that they were mediated by
 231 bacteria, viruses, or apicomplexan pathogens. The numbers of HGTs overlapping with repeats,
 232 including well-known TEs, such as BovB and L1, are shown in the last two columns.

233

Species	Name	HGTs	Novel HGTs	With medium organisms			Overlapped with repeats	With TEs	
				Bacteria	Viruses	Apicomplexa		BovB	L1
anoCar2	Lizard	4	4	1	0	1	1	0	0
bosTau7	Cow	84	74	69	0	43	82	21	56
ce11	Nematode	22	22	1	0	1	0	0	0
danRer10	Zebrafish	25	25	2	0	1	21	4	0
dm6	Fruit fly	177	177	45	1	7	10	0	0
galGal4	Chicken	13	13	1	0	0	1	0	0
hg38	Human	313	152	278	159	268	313	0	117
loxAfr3	Elephant	358	358	148	0	66	317	23	273
mm10	Mouse	15	15	10	0	4	10	0	8
rn6	Rat	13	12	9	0	2	13	0	7
sacCer3	Yeast	27	25	12	0	5	0	0	0
tair10	Arabidopsis	22	22	16	0	3	0	0	0
xenTro9	Frog	17	17	0	0	0	17	0	0

234

235 **Figures**



236

237 **Figure 1. HGTs among eukaryotes.** All 824 eukaryotes were clustered into seven sub-groups:

238 primates, non-primate mammals, non-mammal vertebrates, invertebrates, protozoa, fungi, and

239 plants. (A) The HGT-appearance numbers between 13 model organisms (X axis) and 824

240 eukaryotes (Y axis) are represented by the grid colors in the heatmap. (B) Cross-kingdom HGTs

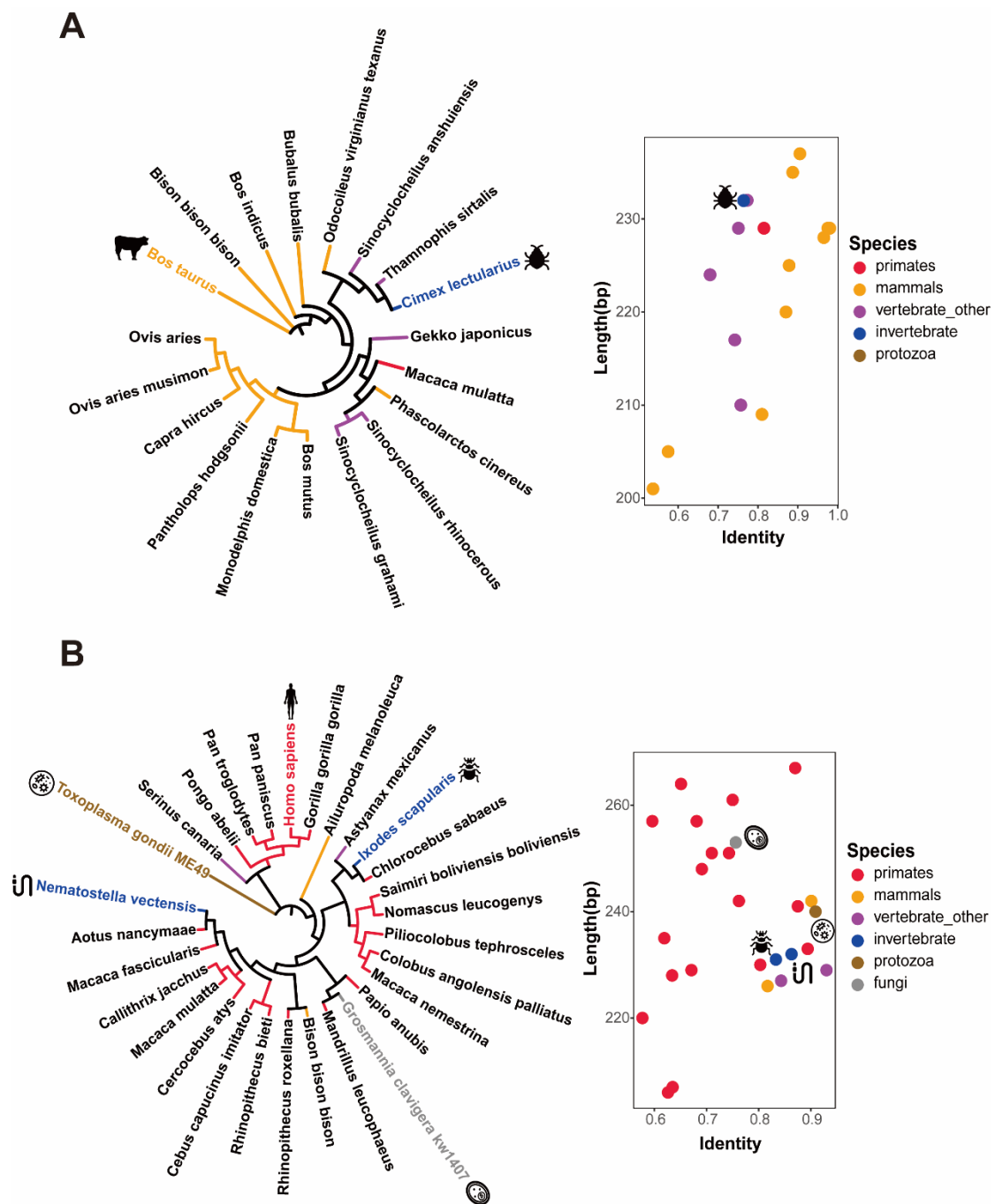
241 were shown by the lines connecting related species, and the thickness of the line represented the

242 HGT-appearance number of the related species. Cross-phylum and cross-class HGTs are shown in

243 Figure S2.

244

245



246

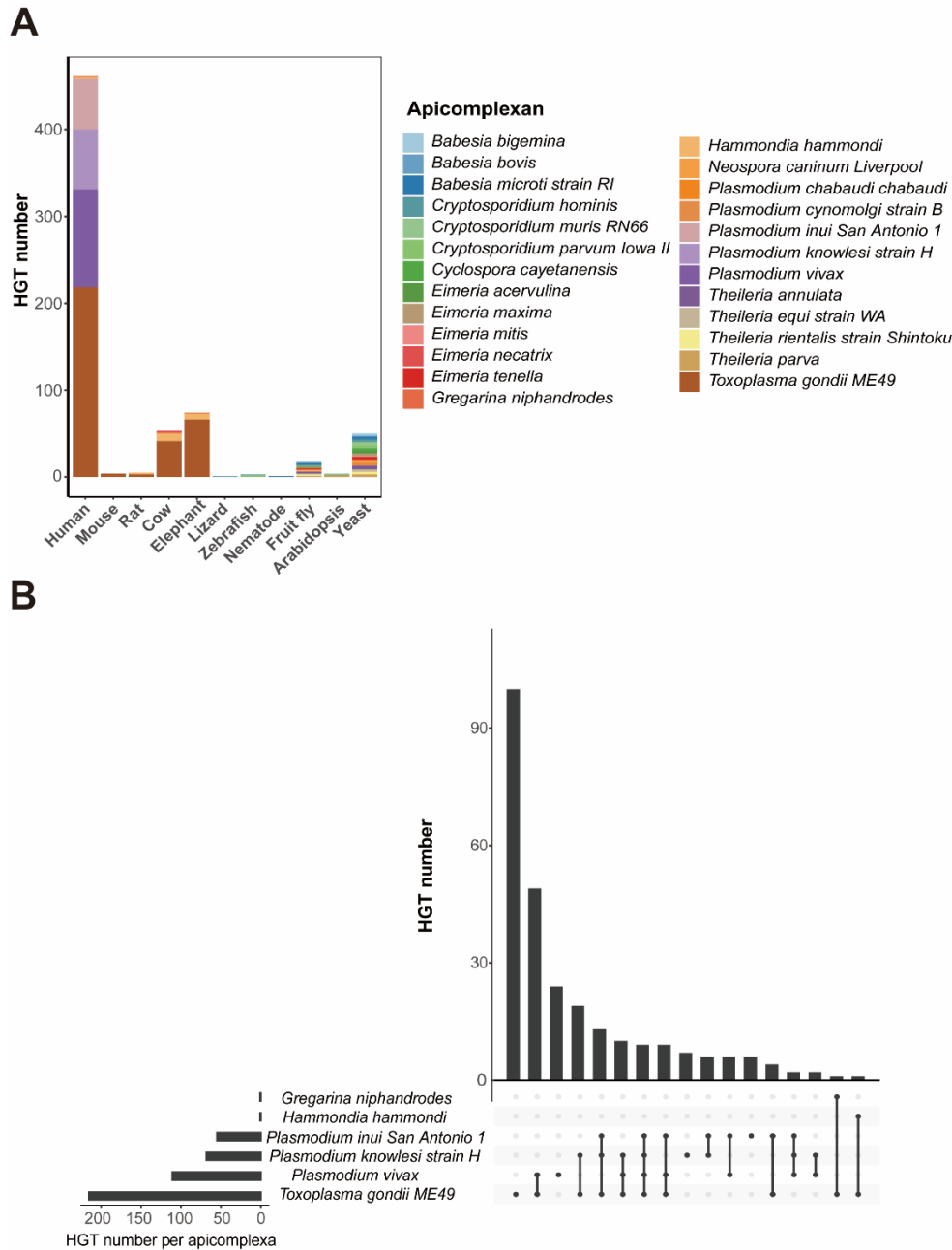
247 **Figure 2. Phylogenetic trees and length-vs-identity plots of HGT region examples. (A) *Bos***

248 *taurus* HGT region “chr25:1343971-1344200”; and (B) *Homo sapiens* HGT region

249 “chr11:24184801-24185043”. The trees on left side represent the evolutionary relationship of

250 species linked by this HGT region, and the plots present sequence similarity between the

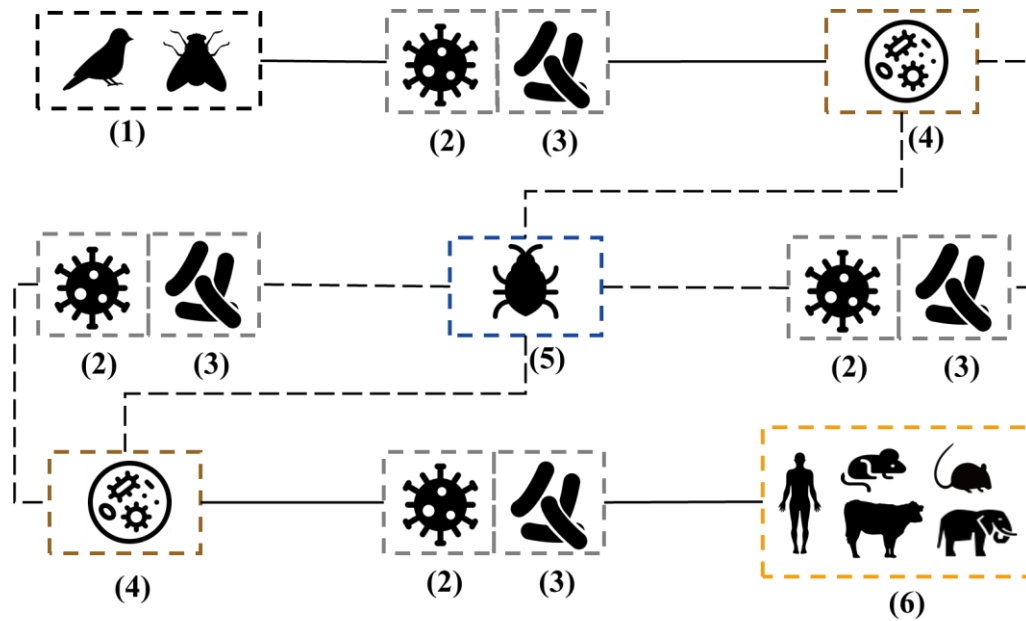
251 homologous sequences from the model organism and the related species.



252

253 **Figure 3. Apicomplexan related HGTs.** (A) The numbers of HGTs associated with
 254 apicomplexans in different model organisms. The X-axis represented 11 different model
 255 organisms and the Y-axis represents the number of corresponding HGTs while different colors
 256 correspond to apicomplexan species. Some HGT sequences from different apicomplexan may
 257 overlap. (B) Detailed information about apicomplexan related HGT regions in human. The X-axis
 258 represented different combination of apicomplexans and the Y-axis represents the numbers of
 259 corresponding HGTs in the human genome.

260



261

262

Figure 4. Putative route of horizontal gene transfer between mammals and distantly related

263

eukaryotes. Here, boxes represent the species that participate in the DNA transfer: (1) distantly

264

related eukaryotes, such as *Drosophila willistoni* and *Serinus canaria*; (2) viral gene pool; (3)

265

bacterial gene pool; (4) intracellular parasites, like *Toxoplasma gondii* ME49; (5) blood-sucking

266

parasites, like *Cimex lectularius*; (6) mammals, including *Homo sapiens*, *Bos taurus*, *Loxodonta*

267

Africana, *Mus musculus*, *Rattus norvegicus*. In this flowchart, solid lines stand for those well

268

supported HGT events in this study, and dashed lines indicate untested hypothesis.

269

270 **Methods**

271 In bacterial genomes, HGT regions are also called genomic islands (GIs) and can be detected using
272 two distinct bioinformatic approaches, based on sequence composition or comparative genomics⁴⁴.
273 In general, the sequence composition of GIs is significantly different from that of the recipient
274 genome. Composition-based methods identify GIs within genome sequences by calculating the k-
275 mer frequencies of a fragment and comparing that frequency distribution with that obtained from
276 the whole genome. Comparative genomics approaches are based on the premise that DNA sequence
277 based phylogenetic tree topology of GIs will be discordant with respect to known species
278 relationships, where sequences that are absent in several closely related organisms appear in more
279 distant species. These two methods can be adapted to the identification of HGTs in eukaryotes but
280 not without challenges. Due to the large sizes and the high heterogeneity of eukaryotic genomes,
281 composition-based approaches may produce a number of false-positive predictions while
282 comparative genomic methods are computationally expensive and time-consuming when hundreds
283 of reference genomes must be aligned. In this study, we identified HGTs between eukaryotes by
284 combining these two approaches to reduce both the false-positive rate and computational cost.

285 **Data collection**

286 Three datasets were downloaded from UCSC Genome Browser⁴⁵
287 (<http://hgdownload.soe.ucsc.edu/downloads.html>) and NCBI Refseq database⁴⁶
288 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq>). The first dataset contained the reference genome
289 sequences of 13 model organisms consisting of 5 mammals, 4 non-mammalian vertebrates, 2
290 invertebrates, 1 fungus and 1 plant. The second dataset, which was used to perform large-scale
291 genomic comparison between model organisms with other species, contained 824 assemblies of
292 eukaryotes including 114 mammals, 125 non-mammalian vertebrates, 155 invertebrates, 81
293 protozoa, 100 plants and 249 fungi. The third dataset, which contained assembled genomes of
294 120,838 bacteria and 7,539 viruses, was created to search the putative intermediary gene pool of
295 DNA during transfer between eukaryotes. Detailed information about these genomes can be found
296 in [Table S11](#).

297 **Pipeline to identify HGTs**

298 [Figure S1](#) shows the pipeline to identify HGTs. Firstly, we identified genomic regions
299 distinguishable from the rest of the genome based on k-mer frequencies (see the next sessions for
300 more details about this genomic fragment filtering step). The selected genomic regions were then
301 aligned with other eukaryotic genomes using LASTZ (version 1.04.00)⁴⁷ ([Supplementary Data 1](#)).
302 A genomic region was considered as a candidate HGT if its sequence level conservation was
303 discordant to its species phylogenetic tree. Specifically, genomic fragments were detected with high
304 identity percentage in species within a distantly related group (DRG) but were missing in most
305 species in a closely related group (CRG). We then clustered the HGT sequences to obtain non-
306 redundant HGTs. Phylogenetic trees for related species and homologous sequences were built for
307 each HGT sequences. related species and homologous sequences of candidate HGT sequences.
308 Finally, each putative HGT region was used to search for homologous sequences in bacteria and
309 viruses; putative agents of HGTs ([Table S6](#)). Detailed information about each step is listed below.

310 **Sequence composition-based genomic fragments filtering**

311 Due to the large sizes of many reference genomes of model organisms, we first screened the
312 potential genomic regions harboring HGT sequences. For each model organism species, we split
313 the genome sequences into 1000-bp segments with 200-bp overlapped regions across all
314 chromosomes. Sequence segments with Ns were left out. Four-bp kmer frequencies were obtained
315 for the whole genome sequences as well as all genome segments. Euclidean distance was used to
316 measure the difference between each segment and the whole genome sequence. All the distances
317 were sorted in descending order. Finally, the fragments whose distances ranked in the top 10% for
318 ce11, dm6, sacCer3 and tair10 due to their smaller genome sizes, top 20% for danRer10 or top 1%
319 for the other model organisms were chosen for further analysis. We tried different kmer sizes (1~6),
320 and k=4 was selected because the highest portion of candidate HGTs previously reported in the
321 human genome¹⁹ were kept ([Figure S6](#)).

322 **Evaluation of the preliminary screening step**

323 Using sequence composition to screen candidate HGT genomic regions was based on the hypothesis
324 that different organisms have different sequence compositions. We tested this hypothesis with
325 available genomes. First, the GC content of whole reference genome sequences showed taxon
326 specific diversity across nine taxonomic clusters ([Figure S7A](#)). Second, principal component

327 analysis (PCA) was performed on 824 eukaryotes using the 4-mer frequencies (Figure S7B). The
328 resulting two-dimensional vectors were then used for binary classification to distinguish whether
329 the organism was a mammal. This approach accurately predicted mammalian genomes 89.35% of
330 the time, with only several non-mammalian vertebrates mis-predicted as mammals (Figure S7C).

331 **PCA and binary-classification**

332 For 824 eukaryotes, we conducted PCA using the “princomp” function in RStudio (version 3.5.0),
333 to reduce the 4-mer frequencies into a lower-dimensional vector. Only the first two principal
334 components, PC1 and PC2, were used as the features to distinguish different species. In the process
335 of binary-classification to determine whether a given species was a mammal, two-dimensional
336 vectors of 824 eukaryotes were randomly divided into a training dataset and a test dataset. The
337 classifier was built from the training dataset using a logistic regression algorithm, which was
338 implemented with the “glm” function in RStudio with default parameters, and the predicted result
339 was evaluated by precision (exactly predicted species/all species in test dataset).

340 **Genome comparison**

341 The genome comparison was conducted using LASTZ for the filtered fragments of the model
342 organisms and the whole genomes of other species with the following arguments: “--format=axt+ -
343 -ambiguous=iupac”.

344 **Re-screening the fragments and search for HGTs**

345 Every organism belongs to its own kingdom, phylum, and class. For each classification level
346 (kingdom, phylum, or class) for each model organism, the other species were separated into two
347 groups: a closely related group (CRG) including all species in the same classification level as the
348 model organism, and a distantly related group (DRG) including all species belonging to different
349 classification levels. For example, when using class as the classification level and using human as
350 the model organism, all mammals were regarded as part of the CRG, while the non-mammalian
351 species formed the DRG. We further screened the filtered fragments based on alignment results
352 from LASTZ, to identify regions with discordant evolutionary relationships. Fragments were
353 regarded as putative HGTs when they had homologs in DRG species but not in the majority of CRG
354 species (see below).

355 The aligned regions (ARs) of the input fragments were retrieved and used to identify putative

356 HGTs. Firstly, we kept ARs that matched to DRG species that were longer than 200bp with a
357 nucleotide identity percentage greater than 70%. For these ARs, we compared the alignment results
358 for CRG species, for which the identity percentage threshold was set to 50%. An AR was considered
359 to be present in a CRG species if it was aligned over 60% of its length. In addition, we counted the
360 frequencies for each AR in CRG species (referred to as “CRG scale”). To reduce false positive
361 results generated by incorrect alignments, we removed ARs that contained the character ‘N’ or
362 whose GC percentages were less than 0.3 or greater than 0.6. Finally, we checked the repetitive
363 regions overlapping with ARs. RepeatMasker tracks were downloaded from the UCSC Genome
364 Browser, or we ran *de novo* RepeatMasker (version 4.0.7)⁴⁸ (<http://www.repeatmasker.org>) to label
365 the repeats of ARs. We then removed ARs that overlapped with simple repeats or low complexity
366 repeats. We also use TRF(version 4.09)⁴⁹ to remove ARs that overlapped with any random repeats.

367 We set M as the maximum number of species with sequences aligned in the CRG. For example,
368 when using class as the classification level and using human as the target model organism, all
369 mammals were regarded as the CRG, while the non-mammalian species formed the DRG. M can
370 be set to the number of all primates which is the order humans belong to. The M values can be set
371 to the numbers of vertebrates, mammals and primates in the genome dataset containing 824
372 eukaryotes (Table S12). At the same time, the DRG scales of ARs were limited such that they
373 appeared in at least N of all the species in the DRG. In our analysis, we set N=1. The alignment
374 threshold for the DRG, including identity and length coverage, were set much higher (see below)
375 than for CRGs, and we removed ARs with high percentages of GC or repeat compositions. The
376 remaining ARs were considered as candidate HGTs, and were used to build trees to determine
377 discordance with known evolutionary relationships. The detailed parameter setting are shown in
378 [Table S12](#).

379 **Identifying non-redundant HGTs**

380 HGTs were clustered using the cd-hit-est program (version 4.6.6)⁵⁰ with minimum nucleotide
381 identity set at 80%. The longest sequences from each cluster were selected to represent the non-
382 redundant HGTs.

383 **Counting the copy numbers of HGTs**

384 We run BLASTN alignment for non-redundant HGT sequences against their host reference

385 genomes, with the parameter “-e 1e-5”. For each HGT, we selected aligned regions that covered at
386 least 90% of HGT regions with nucleotide identity > 90%. We then merged those aligned regions
387 with overlapped coordinates. The copy number of each HGT was determined from the number of
388 merged HGT copies.

389 **Exclusion of mitochondrial or chloroplast DNA**

390 Complete mitochondrial genomes of 13 model organisms and *Arabidopsis thaliana* chloroplast
391 DNA were obtained from NCBI, and we then searched with BLASTN against non-redundant HGTs.
392 With the argument “-evalue 1e-5”, we found no homologous DNA sequences in mitochondrial or
393 chloroplast genomes.

394 **Remove HGTs present in Endogenous viruses**

395 Endogenous retroviruses (ERVs) are widespread in vertebrates, making up nearly 8% of the genome
396 of *Homo sapiens*⁵¹. ERVs in human share sequence homology with other primate ERVs⁵².
397 Therefore, in order to avoid reporting sequences as HGTs that are actually from ancestral
398 inheritance, we removed all HGTs found in ERVs. We collected ERVs from the repeat annotation
399 of the UCSC genome browser, except for *Saccharomyces cerevisiae* S288C and *Arabidopsis*
400 *thaliana*. All HGTs that overlapped with ERV genomic coordinates or aligned to ERVs using
401 BLASTN (identity>90% and length>100bp) were removed.

402 **Comparison with reported HGTs in previous studies**

403 We obtained reported HGTs for these model organisms from previous publications, including
404 genomic coordinates and DNA sequences. HGTs in our study were considered novel if they did not
405 match reported HGTs by genomic coordinates or sequence alignment (BLASTN⁵³, matched
406 length>200bp and identity>80%).

407 **Construction of HGT phylogenetic tree**

408 For each HGT, we searched for homologous sequences in other species based on the LASTZ output.
409 The nucleotide identity threshold of homologous sequences was set at 70% for DRG species and
410 50% for CRG species. When multiple regions in a species met the criteria, the best matched
411 sequence, which had the maximal score weighted by the identity and multiplied by the alignment
412 length, was picked to represent the homologous sequence. Based on the HGT sequence and the
413 homologous sequences collected from other species, we ran multiple sequence alignment using

414 muscle (version 3.8.31)⁵⁴ and then used FastTree (version 2.1.9) to build a maximum likelihood
415 phylogenetic tree, which was visualized with iTOL (version 3.0)⁵⁵. The homologous regions in other
416 species and phylogenetic trees for non-redundant HGTs can be found in [Table S13](#).

417 **Homologous sequences in bacteria and viruses**

418 HGTs were aligned to the assemblies of bacteria and viruses using NCBI BLAST(2.9.0+) with
419 parameters: “-task blastn -evalue 1e-3”. Matched regions in assemblies were filtered to be longer
420 than 200bp, with nucleotide identity greater than 60%.

421 **Validation of homologous sequences in eukaryotic genomes with WGS datasets**

422 Discordant HGT trees, constructed from discordant sequences from reference genomes, were the
423 principal evidence for identifying HGTs from our pipeline. Thus, the power for detecting HGTs
424 depended heavily on the quality of the reference genomes. Contaminating sequences from other
425 species were the most likely sources of false positives. For model organisms, most candidate
426 transferred DNA were also found in their sibling lineages, therefore the probability of sequencing
427 contamination was negligible. However, the inaccurate reference genomes of other eukaryotes (such
428 as parasites and protozoan pathogens) could cause false positive results due to sequencing
429 contamination. For example, if an abnormal HGT tree consists of only one parasite and several
430 primates, and the process of constructing the reference genome of this parasite was contaminated
431 by human DNA, this DNA transfer would be an artifact. We checked for contamination artifacts in
432 candidate transferred DNA by alignment with whole genome sequencing (WGS) raw data from
433 species present in discordant cross-kingdom HGT trees. In total, we collected 59 species which were
434 in different kingdoms with the target model organism, including 15 protozoa, 20 plants, 3 fungi, 11
435 invertebrates, and 10 vertebrates. For each of these species, we downloaded multiple WGS raw
436 datasets (ranging from 3 samples to 201 samples) from the SRA database that were not used to
437 construct the reference genome. In total, we obtained 1,190 WGS samples. Sequence alignment was
438 done using Bowtie2 (version 2.2.4)⁵⁶ with default parameters. For each species, we calculated the
439 length coverage percentage for homologous sequences (M sequences) of WGS samples (N
440 samples), thus generating a coverage percentage matrix (M*N). Once a sequence had coverage of
441 over 80% by any samples, it was classified as not an artifact. The results are shown in [Table S14](#).

442 **Functional annotation of genes influenced by HGTs**

443 Genome annotation files (GFF or GTF format) were obtained for model organisms from Ensembl
444 ⁵⁷(<http://asia.ensembl.org>) and Tair⁵⁸ (<https://www.arabidopsis.org>), and they were used to identify
445 protein-coding genes and non-coding genes likely to be affected by HGTs (overlapping with HGTs
446 with at least 1bp). The Ensembl gene IDs were input to DAVID (version 6.8)⁵⁹
447 (<https://david.ncifcrf.gov>) for functional enrichment analysis. Significantly enriched Gene Ontology
448 terms (GO terms) (Bonferroni<0.05) for these genes were shown in the results.

449 **Evaluation of the pipeline using simulated datasets**

450 We constructed a simulated genome (called genome H) with 175 HGTs from a set of distantly
451 related genomes (called Genome set D) to the human genome. Genome set D has 4 cruciferous plant
452 genomes, including *Arabidopsis thaliana*, *Brassica napus*, *Brassica oleracea var. oleracea* and
453 *Brassica rapa*), while Genome set C contains 4 primate genomes, *Pan paniscus*, *Pan troglodytes*,
454 *Pongo abelii* and *Gorilla gorilla gorilla*. The 175 HGTs are sequences that have high similarity
455 with genomes in Genome set D (>90%) but have low similarity (<10%) with genomes in Genome
456 set C, the closely related group of genomes .

457 Firstly, the genome comparison between genomes in Genome set D was conducted using
458 LASTZ⁴⁷ and Multiz⁶⁰ to obtain sequences whose identity percentages with all genomes of Genome
459 set D were >90% and lengths >200bps. These sequences were compared with the genomes in
460 Genome set C and the sequences having low similarity (identity <10%) were reserved. The obtained
461 sequences were then clustered using the cd-hit-est program (version 4.6.6)⁵⁰ with minimum
462 nucleotide identity set at 80%. The longest sequences from each cluster were selected as simulated
463 HGTs, which were 175 in total. These 175 HGTs were then evenly divided into 10 groups according
464 to their sequence lengths, and their copy numbers were ranged from 2⁰ to 2⁹. Eventually, 175 HGTs
465 with different copy numbers were inserted into the human genome as a simulated genome H.
466 Finally, we ran our pipeline with genome H as the target genome, genome set D as remote genome
467 set, genome set C as closely related genome set and parameters M, N, L as 1, 1, 200 respectively. If
468 the a correct HGT region was covered more than 60% of its length by a predicted HGT region, the
469 prediction was considered correctly predicted.

470

471 **Declarations**

472 **Ethics approval and consent to participate**

473 Not applicable.

474 **Competing interests**

475 The authors declare that they have no competing interests

476 **Authors' contributions**

477 CCW conceived and designed the study. KL, FZY and CCW developed the pipeline and identified

478 HGTs. KL, FZY and ZQD collected the datasets. KL and FZY conducted the visualization. KL,

479 FZY, CCW and DLA wrote the manuscript. KL, FZY, CCW, ZQD and DLA revised the manuscript.

480 All authors read and approved the final manuscript.

481 **Acknowledgements**

482 This work was supported by grants from the National Natural Science Foundation of China

483 (32170643, 61472246 and J1210047), Natural Science Foundation of Shanghai (22ZR1433600 and

484 20ZR1428200), the National Basic Research Program of China (2013CB956103), the National

485 High-Tech R&D Program (863) (2014AA021502), the SJTU JiRLMDS Joint Research Fund

486 (MDS-JF-2019A07), and the Cross-Institute Research Fund of Shanghai Jiao Tong University

487 (YG2017ZD01 and YG2015MS39). The funders had no role in study design, data collection and

488 analysis, decision to publish, or preparation of the manuscript. We thank the High Performance

489 Computing Center at Shanghai Jiao Tong University for the computation.

490 **Data availability**

491 All datasets, supplementary tables and an example of analysis pipeline application are listed in the

492 webpage at <http://cgm.sjtu.edu.cn/hgt> (password: hgt2019passwd) (this webpage will become freely

493 available after this paper is accepted).

494 **Code availability**

495 All scripts used in this study are available in GitHub at <https://github.com/SJTU-CGM/HGT.git>.

496

497 **Supplementary Figures, Tables and Datasets**

498 **There are 7 Figures, 14 Tables and 1 datasets provided in multiple supplementary files.**

499 **Descriptions about the figures, tables and datasets are listed below. Supplementary figures**

500 **are listed at the end of this file, while supplementary tables and datasets are accessible from**

501 **the given URL listed in the data availability.**

502

503 **Supplementary Figure 1**

504 The HGT identification system for model eukaryotes.

505 **Supplementary Figure 2**

506 Cross-phylum HGTs and cross-class HGTs.

507 **Supplementary Figure 3**

508 Repeat characteristics of HGT regions as well as reference genomes.

509 **Supplementary Figure 4**

510 The number of HGTs associated with apicomplexan in different model organisms.

511 **Supplementary Figure 5.**

512 Phylogenetic trees of other HGT region examples.

513 **Supplementary Figure 6**

514 The impact of parameter setting for the fast HGT selection step using k-mer frequency. The

515 parameters are k-mer size and fragment percentage.

516 **Supplementary Figure 7**

517 Evaluation of the preliminary screening step.

518

519 **Supplementary Table 1**

520 Detailed information of non-redundant HGTs after removing redundancy, including genomic

521 coordinates, CRG scale, bacterial presence, viral presence, and their copy numbers in the whole

522 genomes.

523 **Supplementary Table 2**

524 HGT-appearance numbers between the 13 model organisms and 824 eukaryotes.

525 **Supplementary Table 3**

526 The number of cross-kingdom HGTs, cross-phylum HGTs and cross-class HGTs.

527 **Supplementary Table 4**

528 The character of media organisms of HGT regions overlapped with BovB in bosTau7.

529 **Supplementary Table 5**

530 Examples of HGTs in cow and human genomes.

531 **Supplementary Table 6**

532 BLASTN results of non-redundant HGTs against bacteria/viruses.

533 **Supplementary Table 7**

534 Coverage matrices of WGS data for HGT homologous sequences in selected eukaryotes.

535 **Supplementary Table 8**

536 The geographic information of species with HGTs in mammals.

537 **Supplementary Table 9**

538 Putative media of horizontal gene transfer between mammals and distantly related
539 eukaryotes.

540 **Supplementary Table 10**

541 Functional annotation for genes affected by HGTs.

542 **Supplementary Table 11**

543 Information of 13 model organisms and assembly ID of other eukaryotes, bacteria, and viruses.

544 **Supplementary Table 12**

545 Parameter settings of the HGT identification pipeline.

546 **Supplementary Table 13**

547 Homologous regions in other species and phylogenetic trees for non-redundant HGTs.

548 **Supplementary Table 14**

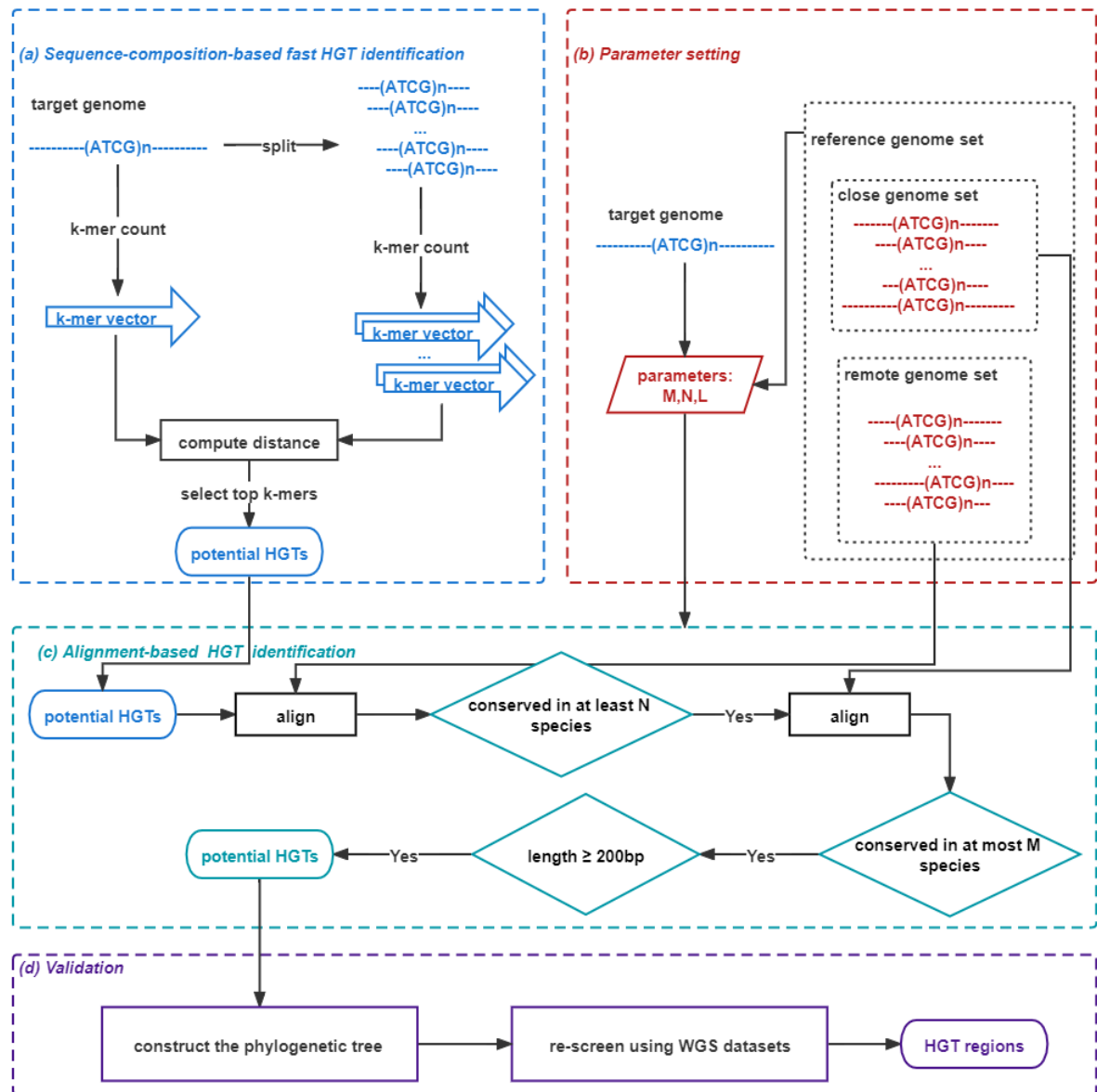
549 Coverage matrices of WGS data for HGT homologous sequences in selected apicomplexan.

550

551 **Supplementary Data 1**

552 Raw output of LASTZ alignment between 13 model organisms with other eukaryotes (197GB)

553 URL: http://cgm.sjtu.edu.cn/hgt/data/Supplementary_Data_1.tar



554

555 **Supplementary Figure 1. The HGT identification system for model eukaryotes.** Four main

556 processes were involved: (a) Sequence-composition-based fast HGT identification: split the

557 chromosomes into fixed sized fragments and screen the fragments to get genome regions with the

558 potential to harbor HGTs according to k-mer frequencies; (b) Parameter setting: choose appropriate

559 parameters according to the target genome and the reference genome sets; (c) Alignment-based

560 HGT identification: sequence comparison between potential HGT containing fragments with the

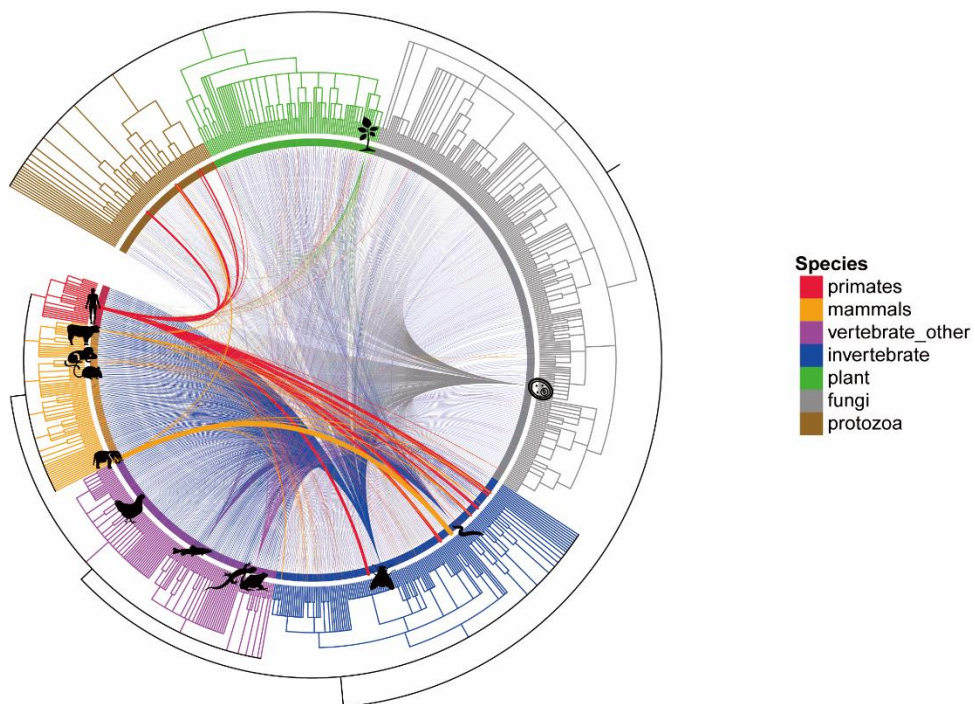
561 whole genomes of other eukaryotes; (d) Validation: re-screen the fragments based on the differences

562 between the HGT region phylogenic tree and the organism phylogenetic tree and further screen the

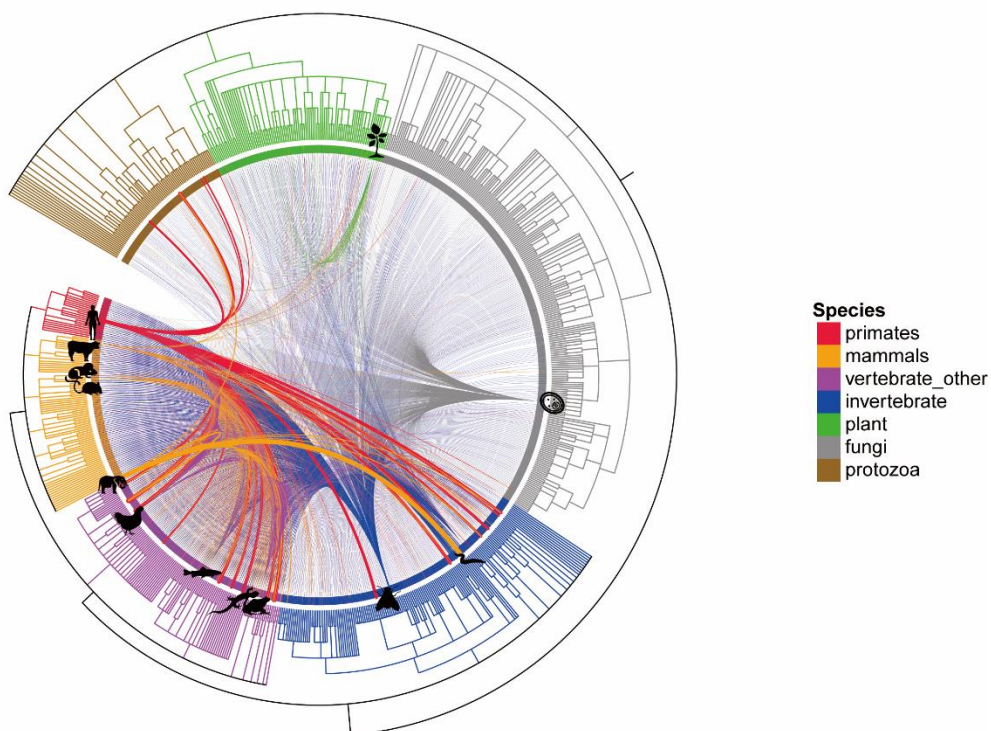
563 fragments using WGS datasets of selected organisms.

564

A



B



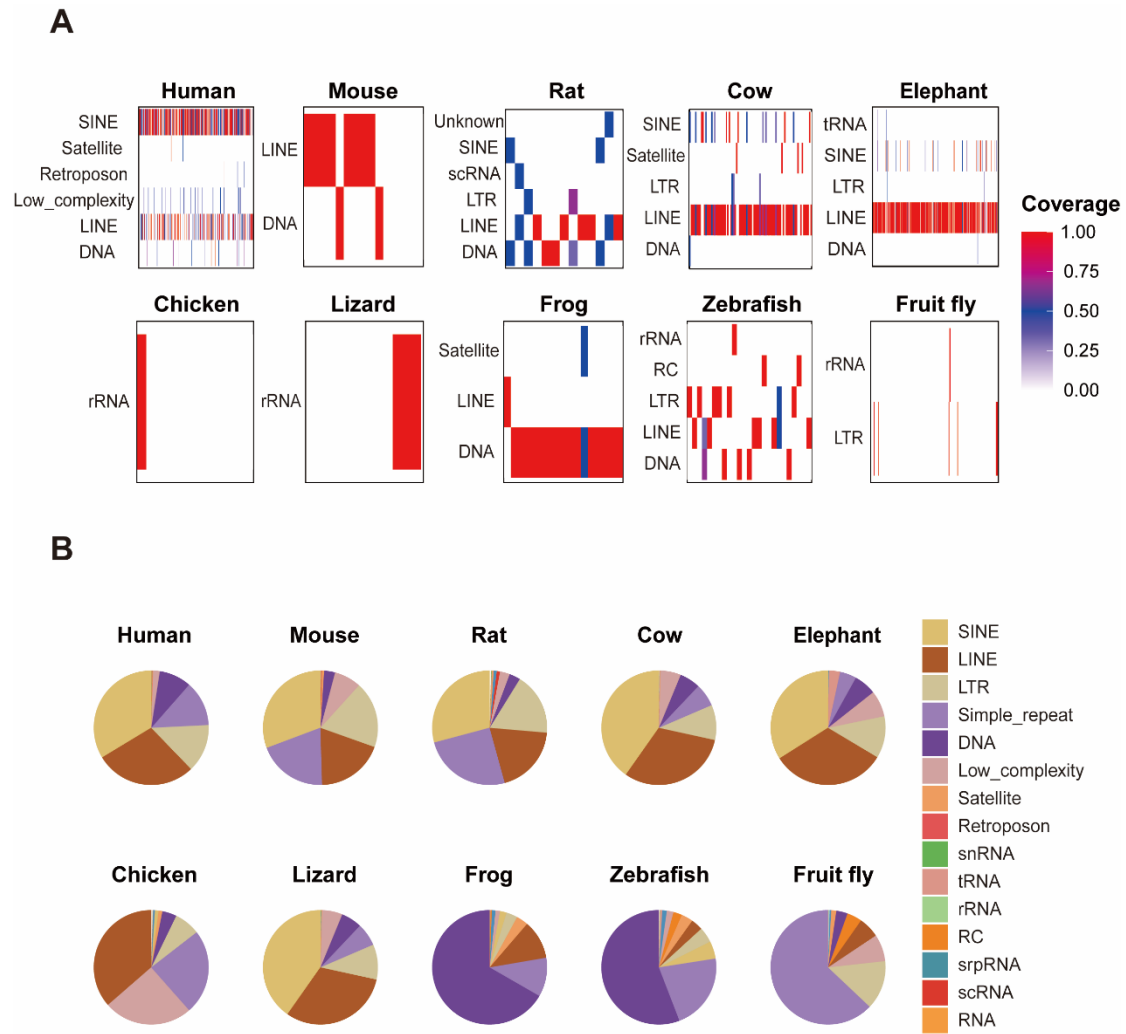
565

566 **Supplementary Figure 2. Cross-phylum HGTs and cross-class HGTs. (A)** Cross-phylum HGTs

567 and **(B)** cross-class HGTs were shown by the lines connecting related species, and the thickness of

568 the line represented the HGT-appearance number of the related species.

569



570

571 **Supplementary Figure 3. Repetitive regions overlapping HGT regions.** (A) Repetitive

572 composition of HGT regions. X-axis represents HGTs ranked by CRG scale in ascending order, and

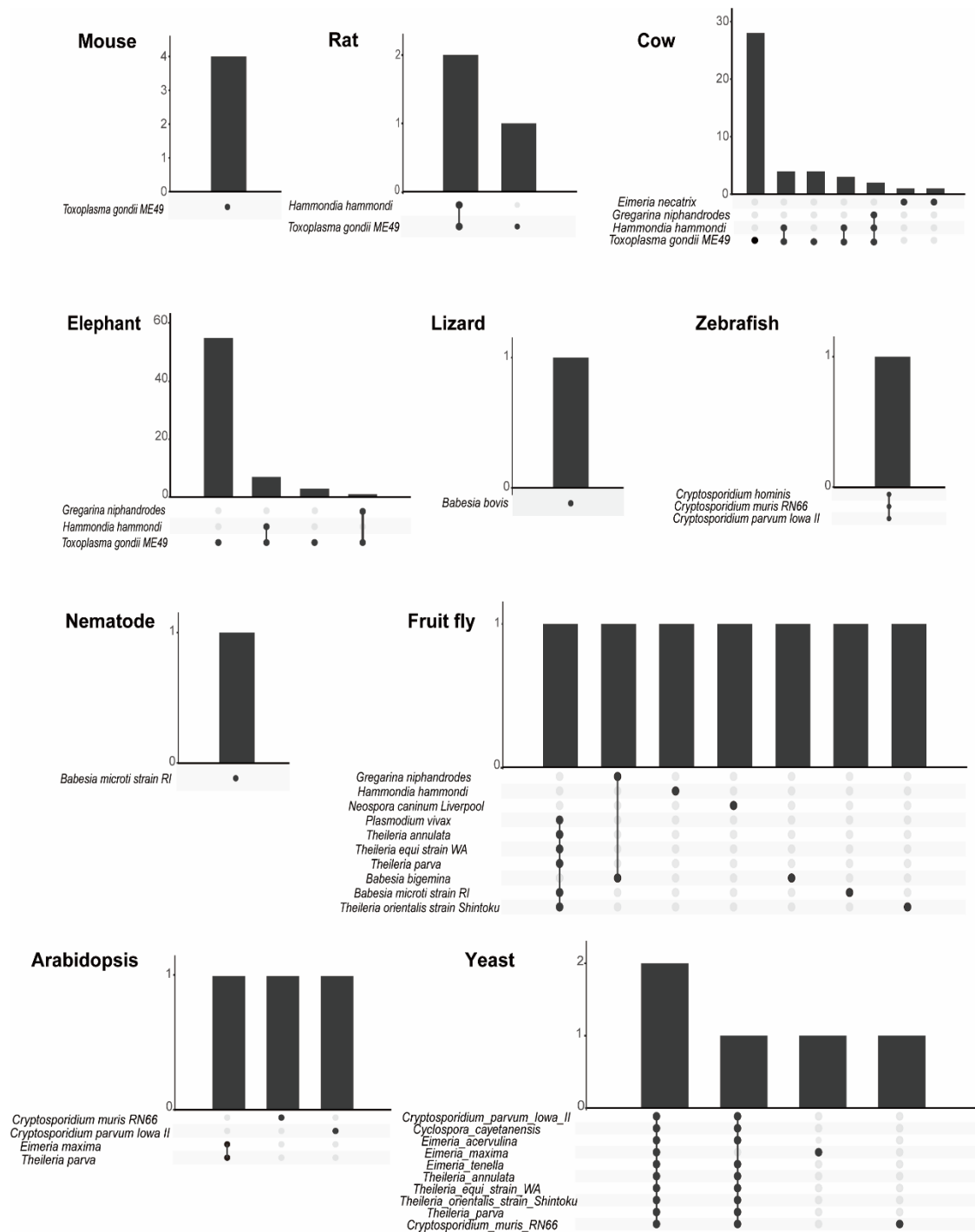
573 the color reflects the percentage of the length of an HGT annotated as a corresponding type of repeat.

574 If a repeat type was not annotated within a given HGT, the color was white; (B) Repeat type

575 composition of the 10 model genomes.

576

577



578

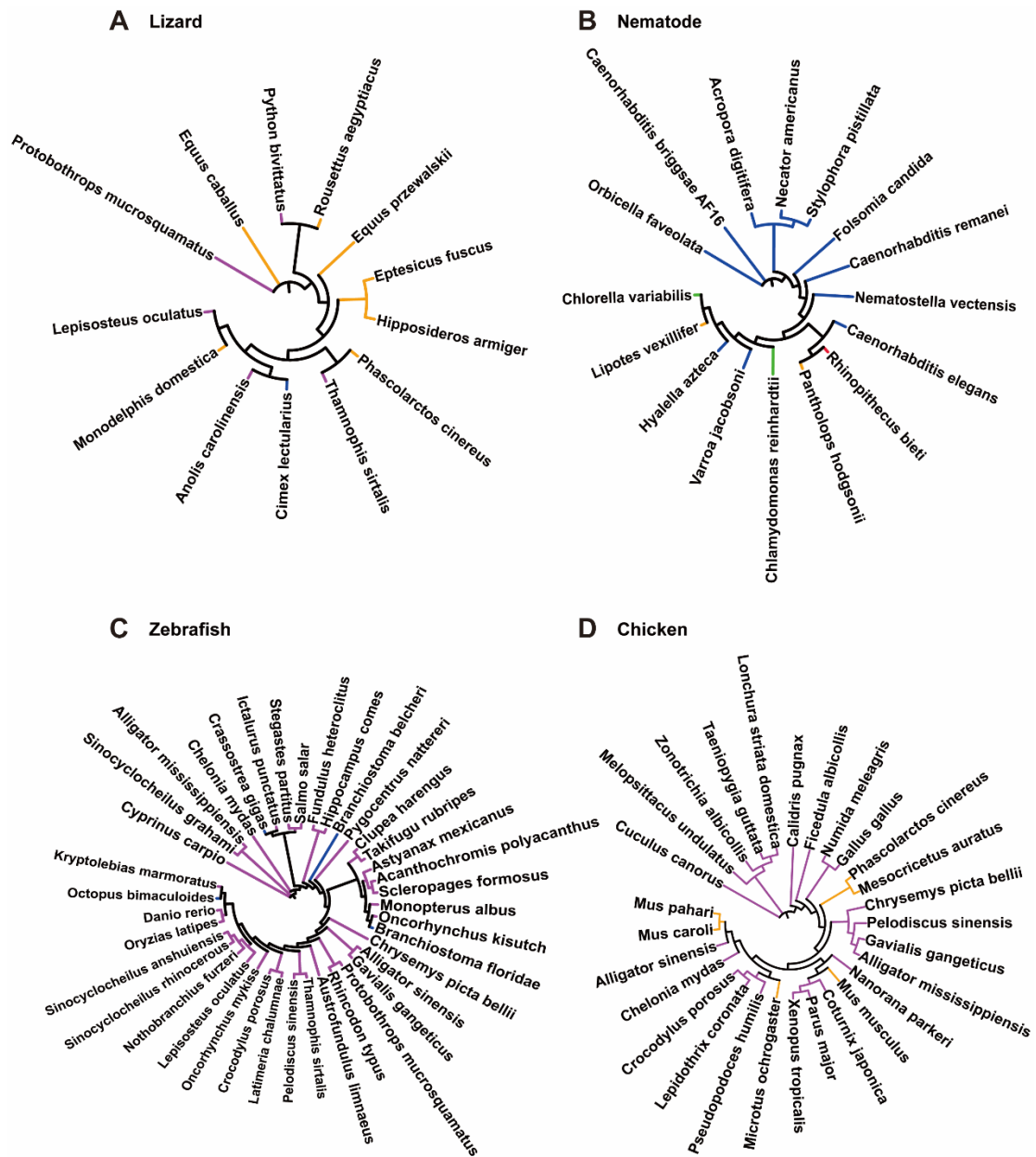
579 **Supplementary Figure 4. The number of HGTs associated with Apicomplexan in different**

580 **model organisms.** For each picture, the X-axis represented different combination of diverse

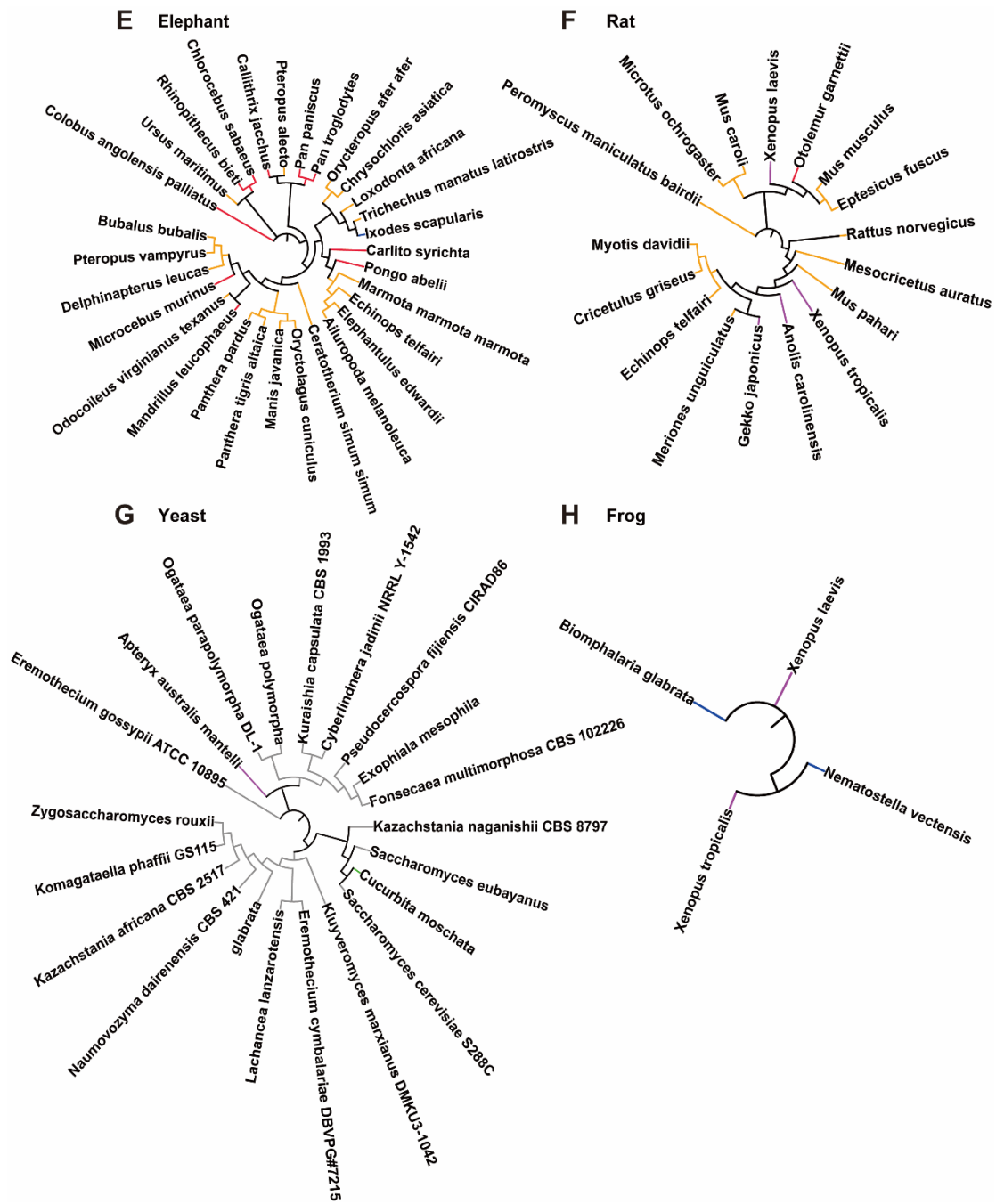
581 apicomplexans and the Y-axis represents the number of corresponding HGTs of the model

582 organism.

583

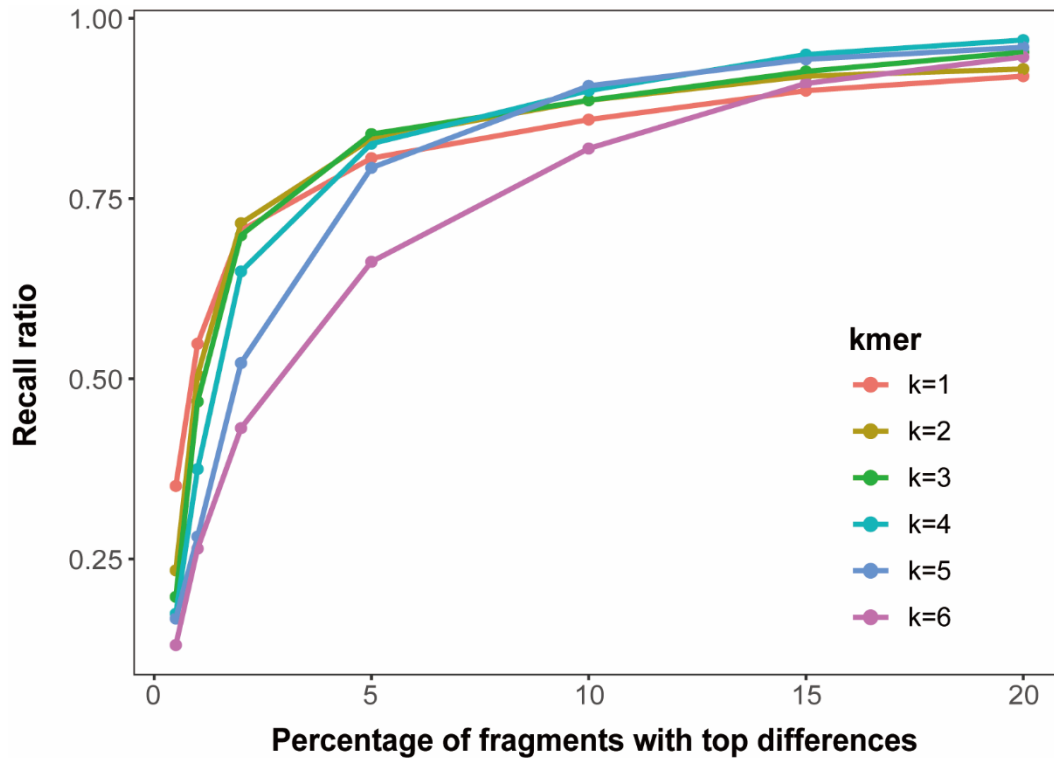


584



585

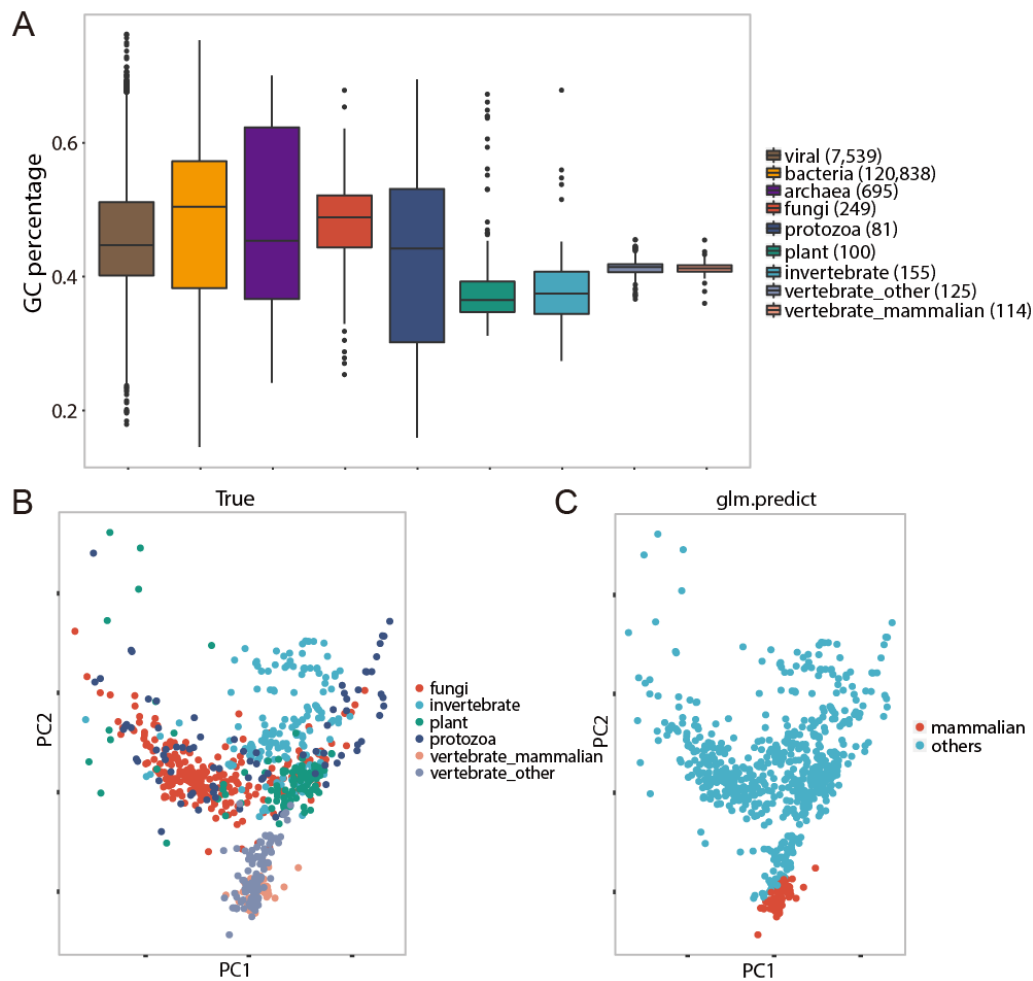
598



599

600

601 **Supplementary Figure 6. The impact of parameter setting on the HGT recalling rate. The**
602 **parameters are k-mer size and fragments percentage.** The X-axis represented the percentage of
603 retained short sequences (0.5%, 1%, 2%, 5% ,10%,15% and 20% were tested), and the Y-axis
604 represented the HGT Recall rate (the number of HGT correctly predicted/ the total number of HGT
605 reported in a previous study). Here only 299 original reported human HGT regions remained after
606 we had filtered out the potential HGT regions containing more than 50% of simple repeats identified
607 by TRF. When k=4, more than 75% of the 299 human HGT regions were identified by our pipeline
608 if only the 5% of the fragments with highest k-mer frequency distances to the human genome were
609 input to the second stage of our pipeline.

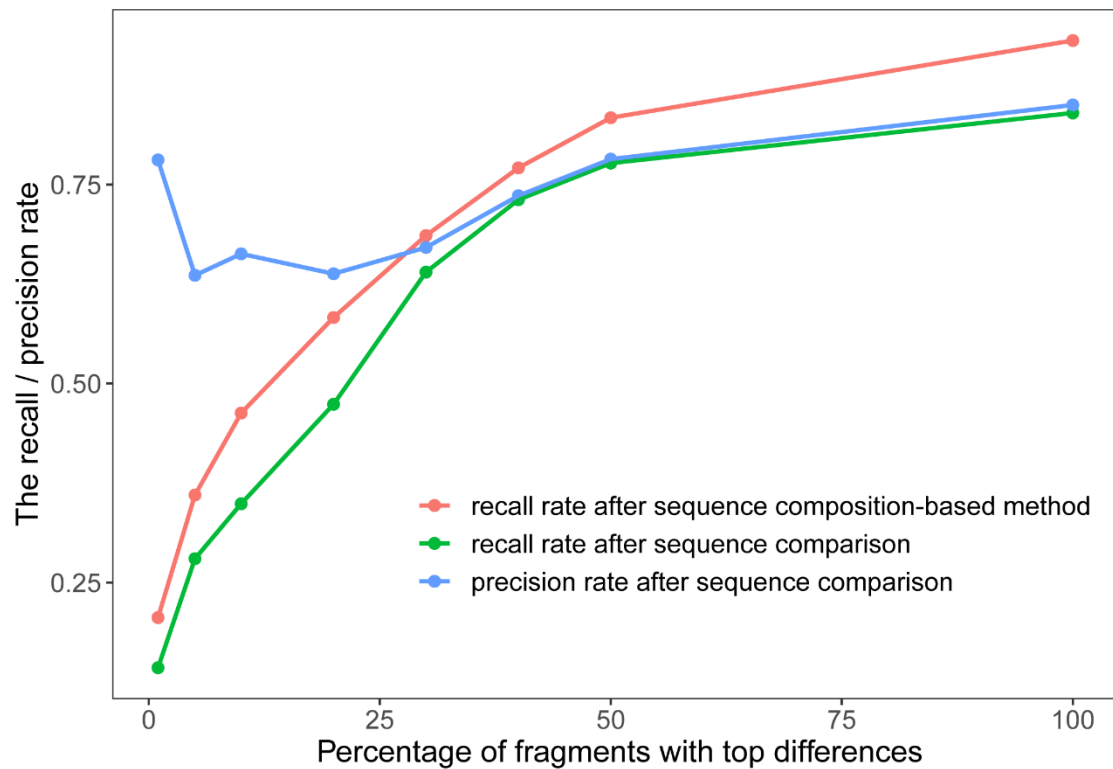


610

611 **Supplementary Figure 7. Sequence composition versus genome classification.** (A) GC
612 percentage distribution of the nine taxonomic categories of organisms; (B) The distribution of the
613 two-dimensional vector (PC1 and PC2) of 824 eukaryotes in six categories; (C) Binary-
614 classification based on the two-dimensional feature, in which the PC1 and PC2 coordinates remain
615 unchanged but the labels of the organism differ according to the predicted results.

616

617



618

619 **Supplementary Figure 8. Evaluation results of the HGT identification pipeline on the**
620 **simulated dataset.** The X-axis represented the percentage of genomic fragments past the sequence
621 composition filtering step (only the top x% of k-mers with largest differences from the genome
622 entered the genome comparison step, where top x% stood for 1%, 5% ,10%,20%, 30%, 40%, 50%
623 and 100%), and the Y-axis represented the HGT recall rate (the number of HGTs correctly predicted/
624 the total number of HGTs in the simulated dataset) or precision rate (the number of HGT correctly
625 predicted/ the total number of HGT predicted).

626

627

628 **References**

- 629 1 Soucy, S. M., Huang, J. & Gogarten, J. P. Horizontal gene transfer: building the web of life.
630 *Nature reviews. Genetics* **16**, 472-482, doi:10.1038/nrg3962 (2015).
- 631 2 Polz, M. F., Alm, E. J. & Hanage, W. P. Horizontal gene transfer and the evolution of
632 bacterial and archaeal population structure. *Trends in genetics : TIG* **29**, 170-175,
633 doi:10.1016/j.tig.2012.12.006 (2013).
- 634 3 Dagan, T., Artzy-Randrup, Y. & Martin, W. Modular networks and cumulative impact of
635 lateral transfer in prokaryote genome evolution. *Proceedings of the National Academy of*
636 *Sciences of the United States of America* **105**, 10039-10044,
637 doi:10.1073/pnas.0800679105 (2008).
- 638 4 Cheng, S. F. *et al.* Genomes of Subaerial Zygnematophyceae Provide Insights into Land
639 Plant Evolution. *Cell* **179**, 1057-+, doi:10.1016/j.cell.2019.10.019 (2019).
- 640 5 Xia, J. X. *et al.* Whitefly hijacks a plant detoxification gene that neutralizes plant toxins. *Cell*
641 **184**, 1693-+, doi:10.1016/j.cell.2021.02.014 (2021).
- 642 6 Leclercq, S. *et al.* Birth of a W sex chromosome by horizontal transfer of Wolbachia
643 bacterial symbiont genome. *P Natl Acad Sci USA* **113**, 15036-15041,
644 doi:10.1073/pnas.1608979113 (2016).
- 645 7 Kado, T. & Innan, H. Horizontal Gene Transfer in Five Parasite Plant Species in
646 Orobanchaceae. *Genome Biol Evol* **10**, 3196-3210, doi:10.1093/gbe/evy219 (2018).
- 647 8 Lukes, J. & Husnik, F. Microsporidia: A Single Horizontal Gene Transfer Drives a Great Leap
648 Forward. *Curr Biol* **28**, R712-R715, doi:10.1016/j.cub.2018.05.031 (2018).
- 649 9 Gilbert, C., Schaack, S., Pace, J. K., Brindley, P. J. & Feschotte, C. A role for host-parasite
650 interactions in the horizontal transfer of transposons across phyla. *Nature* **464**, 1347-
651 U1344, doi:10.1038/nature08939 (2010).
- 652 10 Walsh, A. M., Kortschak, R. D., Gardner, M. G., Bertozzi, T. & Adelson, D. L. Widespread
653 horizontal transfer of retrotransposons. *P Natl Acad Sci USA* **110**, 1012-1016,
654 doi:10.1073/pnas.1205856110 (2013).
- 655 11 Ivancevic, A. M., Kortschak, R. D., Bertozzi, T. & Adelson, D. L. Horizontal transfer of BovB
656 and L1 retrotransposons in eukaryotes. *Genome Biol* **19**, 85, doi:10.1186/s13059-018-
657 1456-7 (2018).
- 658 12 Huang, J. L. Horizontal gene transfer in eukaryotes: The weak-link model. *Bioessays* **35**,
659 868-875, doi:10.1002/bies.201300007 (2013).
- 660 13 Martin, W. F. Too Much Eukaryote LGT. *Bioessays* **39**, doi:10.1002/bies.201700115 (2017).
- 661 14 Salzberg, S. L. Horizontal gene transfer is not a hallmark of the human genome. *Genome*
662 *Biol* **18**, 85, doi:10.1186/s13059-017-1214-2 (2017).
- 663 15 Leger, M. M., Eme, L., Stairs, C. W. & Roger, A. J. Demystifying Eukaryote Lateral Gene
664 Transfer. *Bioessays* **40**, doi:10.1002/bies.201700242 (2018).
- 665 16 Keeling, P. J. & Palmer, J. D. Horizontal gene transfer in eukaryotic evolution. *Nat Rev*
666 *Genet* **9**, 605-618, doi:10.1038/nrg2386 (2008).
- 667 17 Xia, J. *et al.* Whitefly hijacks a plant detoxification gene that neutralizes plant toxins. *Cell*
668 **184**, 3588, doi:10.1016/j.cell.2021.06.010 (2021).
- 669 18 Ivancevic, A. M., Kortschak, R. D., Bertozzi, T. & Adelson, D. L. Horizontal transfer of BovB
670 and L1 retrotransposons in eukaryotes. *Genome Biol* **19**, doi:10.1186/s13059-018-1456-
671 7 (2018).
- 672 19 Huang, W. *et al.* Widespread of horizontal gene transfer in the human genome. *BMC*
673 *Genomics* **18**, 274, doi:10.1186/s12864-017-3649-y (2017).
- 674 20 Keeling, P. J. & Palmer, J. D. Lateral transfer at the gene and subgenomic levels in the
675 evolution of eukaryotic enolase. *P Natl Acad Sci USA* **98**, 10745-10750, doi:DOI
676 10.1073/pnas.191337098 (2001).
- 677 21 Lang, J. L., Gonzalez-Mula, A., Taconnat, L., Clement, G. & Faure, D. The plant GABA

- 678 signaling downregulates horizontal transfer of the *Agrobacterium tumefaciens* virulence
679 plasmid. *New Phytol* **210**, 974-983, doi:10.1111/nph.13813 (2016).
- 680 22 Ge, Y. L. *et al.* Gene transfer of the *Caenorhabditis elegans* n-3 fatty acid desaturase
681 inhibits neuronal apoptosis. *J Neurochem* **82**, 1360-1366, doi:DOI 10.1046/j.1471-
682 4159.2002.01077.x (2002).
- 683 23 Wu, B. *et al.* Interdomain lateral gene transfer of an essential ferrochelatase gene in human
684 parasitic nematodes. *P Natl Acad Sci USA* **110**, 7748-7753, doi:10.1073/pnas.1304049110
685 (2013).
- 686 24 Sun, B. F. *et al.* Horizontal functional gene transfer from bacteria to fishes. *Sci Rep-Uk* **5**,
687 doi:10.1038/srep18676 (2015).
- 688 25 Brown, A. N. & Lloyd, V. K. Evidence for horizontal transfer of *Wolbachia* by a *Drosophila*
689 mite. *Exp Appl Acarol* **66**, 301-311, doi:10.1007/s10493-015-9918-z (2015).
- 690 26 Palazzo, A., Lovero, D., D'Addabbo, P., Caizzi, R. & Marsano, R. M. Identification of Bari
691 Transposons in 23 Sequenced *Drosophila* Genomes Reveals Novel Structural Variants,
692 MITEs and Horizontal Transfer. *Plos One* **11**, doi:10.1371/journal.pone.0156014 (2016).
- 693 27 Bartolome, C., Bello, X. & Maside, X. Widespread evidence for horizontal transfer of
694 transposable elements across *Drosophila* genomes. *Genome Biol* **10**, doi:10.1186/gb-
695 2009-10-2-r22 (2009).
- 696 28 Pace, J. K., Gilbert, C., Clark, M. S. & Feschotte, C. Repeated horizontal transfer of a DNA
697 transposon in mammals and other tetrapods. *P Natl Acad Sci USA* **105**, 17023-17028,
698 doi:10.1073/pnas.0806548105 (2008).
- 699 29 Crisp, A., Boschetti, C., Perry, M., Tunnacliffe, A. & Micklem, G. Expression of multiple
700 horizontally acquired genes is a hallmark of both vertebrate and invertebrate genomes.
701 *Genome Biol* **16**, doi:10.1186/s13059-015-0607-3 (2015).
- 702 30 Carr, M., Bensasson, D. & Bergman, C. M. Evolutionary Genomics of Transposable
703 Elements in *Saccharomyces cerevisiae*. *Plos One* **7**, doi:10.1371/journal.pone.0050978
704 (2012).
- 705 31 Hall, C. & Dietrich, F. S. The reacquisition of biotin prototrophy in *Saccharomyces*
706 *cerevisiae* involved horizontal gene transfer, gene duplication and gene clustering.
707 *Genetics* **177**, 2293-2307, doi:DOI 10.1534/genetics.107.074963 (2007).
- 708 32 Novick, P., Smith, J., Ray, D. & Boissinot, S. Independent and parallel lateral transfer of
709 DNA transposons in tetrapod genomes. *Gene* **449**, 85-94, doi:10.1016/j.gene.2009.08.017
710 (2010).
- 711 33 Jurka, J., Kapitonov, V. V., Kohany, O. & Jurka, M. V. Repetitive sequences in complex
712 genomes: Structure and evolution. *Annu Rev Genom Hum G* **8**, 241-259,
713 doi:10.1146/annurev.genom.8.080706.092416 (2007).
- 714 34 Doggett, S. L., Dwyer, D. E., Penas, P. F. & Russell, R. C. Bed Bugs: Clinical Relevance and
715 Control Options. *Clin Microbiol Rev* **25**, 164-+, doi:10.1128/Cmr.05015-11 (2012).
- 716 35 Goddard, J. & deShazo, R. Bed Bugs (*Cimex lectularius*) and Clinical Consequences of Their
717 Bites. *Jama-J Am Med Assoc* **301**, 1358-1366, doi:DOI 10.1001/jama.2009.405 (2009).
- 718 36 Alexander, W. G., Wisecaver, J. H., Rokas, A. & Hittinger, C. T. Horizontally acquired genes
719 in early-diverging pathogenic fungi enable the use of host nucleosides and nucleotides.
720 *P Natl Acad Sci USA* **113**, 4116-4121, doi:10.1073/pnas.1517242113 (2016).
- 721 37 Kim, K. & Weiss, L. M. *Toxoplasma gondii*: the model apicomplexan. *Int J Parasitol* **34**,
722 423-432, doi:10.1016/j.ijpara.2003.12.009 (2004).
- 723 38 van Helden, P. D., van Helden, L. S. & Hoal, E. G. One world, one health. *Embo Rep* **14**,
724 497-501, doi:10.1038/embor.2013.61 (2013).
- 725 39 Montoya, J. G. & Liesenfeld, O. Toxoplasmosis. *Lancet* **363**, 1965-1976, doi:Doi
726 10.1016/S0140-6736(04)16412-X (2004).
- 727 40 Alsmark, C. *et al.* Patterns of prokaryotic lateral gene transfers affecting parasitic microbial
728 eukaryotes. *Genome Biol* **14**, doi:10.1186/gb-2013-14-2-r19 (2013).

- 729 41 Seton, M. *et al.* Global continental and ocean basin reconstructions since 200 Ma. *Earth-*
730 *Sci Rev* **113**, 212-270, doi:10.1016/j.earscirev.2012.03.002 (2012).
- 731 42 Goswami, A. A dating success story: genomes and fossils converge on placental mammal
732 origins. *Evodevo* **3**, doi:10.1186/2041-9139-3-18 (2012).
- 733 43 Gheerbrant, E. Paleocene emergence of elephant relatives and the rapid radiation of
734 African ungulates. *Proc Natl Acad Sci U S A* **106**, 10717-10721,
735 doi:10.1073/pnas.0900251106 (2009).
- 736 44 Langille, M. G. I., Hsiao, W. W. L. & Brinkman, F. S. L. Detecting genomic islands using
737 bioinformatics approaches. *Nat Rev Microbiol* **8**, 372-382, doi:10.1038/nrmicro2350
738 (2010).
- 739 45 Casper, J. *et al.* The UCSC Genome Browser database: 2018 update. *Nucleic Acids Res* **46**,
740 D762-D769, doi:10.1093/nar/gkx1020 (2018).
- 741 46 Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a curated
742 non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids*
743 *Res* **35**, D61-D65, doi:10.1093/nar/gkl842 (2007).
- 744 47 Harris, R. S. Improved pairwise alignment of genomic dna. (2007).
- 745 48 Price, A. L., Jones, N. C. & Pevzner, P. A. De novo identification of repeat families in large
746 genomes. *Bioinformatics* **21**, I351-I358, doi:10.1093/bioinformatics/bti1018 (2005).
- 747 49 Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids*
748 *Res* **27**, 573-580, doi:DOI 10.1093/nar/27.2.573 (1999).
- 749 50 Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of
750 protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659,
751 doi:10.1093/bioinformatics/btl158 (2006).
- 752 51 Paces, J., Pavlicek, A. & Paces, V. HERVd: database of human endogenous retroviruses.
753 *Nucleic Acids Res* **30**, 205-206, doi:DOI 10.1093/nar/30.1.205 (2002).
- 754 52 Johnson, W. E. Endogenous Retroviruses in the Genomics Era. *Annu Rev Virol* **2**, 135-159,
755 doi:10.1146/annurev-virology-100114-054945 (2015).
- 756 53 Camacho, C. *et al.* BLAST plus : architecture and applications. *Bmc Bioinformatics* **10**,
757 doi:10.1186/1471-2105-10-421 (2009).
- 758 54 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
759 throughput. *Nucleic Acids Res* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).
- 760 55 Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and
761 annotation of phylogenetic and other trees. *Nucleic Acids Res* **44**, W242-W245,
762 doi:10.1093/nar/gkw290 (2016).
- 763 56 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
764 **9**, 357-U354, doi:10.1038/Nmeth.1923 (2012).
- 765 57 Zerbino, D. R. *et al.* Ensembl 2018. *Nucleic Acids Res* **46**, D754-D761,
766 doi:10.1093/nar/gkx1098 (2018).
- 767 58 Poole, R. L. The TAIR database. *Methods Mol Biol* **406**, 179-212, doi:10.1007/978-1-
768 59745-535-0_8 (2007).
- 769 59 Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large
770 gene lists using DAVID bioinformatics resources. *Nat Protoc* **4**, 44-57,
771 doi:10.1038/nprot.2008.211 (2009).
- 772 60 Blanchette, M. *et al.* Aligning multiple genomic sequences with the threaded blockset
773 aligner. *Genome Res* **14**, 708-715, doi:10.1101/gr.1933104 (2004).
- 774